

Adaptation of fictional and online conversations to communication media

Christian M. Alis and May T. Lim

National Institute of Physics,
University of the Philippines Diliman
1101 Quezon City, Philippines

Received: date / Revised version: date

Abstract. Conversations allow the quick transfer of short bits of information and it is reasonable to expect that changes in communication medium affect how we converse. Using conversations in works of fiction and in an online social networking platform, we show that the utterance length of conversations is slowly shortening with time but adapts more strongly to the constraints of the communication medium. This indicates that the introduction of any new medium of communication can affect the way natural language evolves.

PACS. 89.65.Ef Social organizations; anthropology – 89.20.-a Interdisciplinary applications of physics

1 Introduction

With an estimated vocabulary size of 20,000 to 40,000 base words [1,2,3], conversations quickly transfer short bits of information via two general means: the oral and the written form. Although the written vocabulary is often larger [4], the grammatically looser and more error-prone oral medium has the advantage of having access to nonverbal cues like gestures and intonations [5] to aid communication. Aside from vocabulary size—word choices, unconsciously repeating words, and other idiosyncrasies [6] also affect the way we perceive conversations.

Conversation analysis typically looks into how turn taking patterns in institutional settings depart from those observed in informal conversations [7], or on the psychological or sociological aspects [8] of social structure. In this work, the length distribution of a single speaking turn, or utterance, was derived to determine if the medium affects the way we express ideas by using datasets that include a mix of real-world (online) and fictional (offline) conversations: online conversation in Twitter (twitter.com); conversations from 19th century novels and short stories; and subtitles from 20th century movies.

Humans typically converse orally, thus the analysis of conversations is usually performed by transcribing recorded audio conversations into text. In cases when this is not possible e.g., before the invention of recorded audio, one technique is to use written records of real and constructed conversations as were done in studies on the emergence of complementary clauses (Paul persuaded John *to kiss Mary*) [9], the use of *do* in negative declaratives (I *do not* understand you) [10], and the increasing prevalence of the modals *gonna*, *gotta* and *wanna* [11]. Written records of

spoken speech are also included in corpora like *A Corpus of English Dialogues 1560–1760* [12] and *The Corpus of Historical American English: 400 million words, 1810–2009* [13]. However, only conversations (fictional dialogues) in novels, short stories, and movies were analyzed in this paper because utterances tend to be less narrative and directed to another person unlike in other genres like drama comedies or trial transcripts. Although it has been shown that styles vary across and even within authors [14], we assumed that conversations in their works are mostly independent of the author's style, i.e., a conversation in their works conveys how another person (character), and not how the author, speaks. Furthermore, errors due to transcribing are practically eliminated when using books and movies.

Twitter, as a form of computer-mediated communication, is different from oral or written media [15]. While assumed to be happening in real-time, the purely written nature of a Twitter-based conversation differentiates it from the transcribed oral communication in books and movies. In addition, Twitter conversations have an explicit length limit—an utterance can only be up to 140 characters long.

Putting a length constraint on the outset would show drastic changes. A case in point would be SMS messages. At its peak, *textspeak* looked very much different from standard spelling—primarily due to the effort it takes to spell out words through a numerical keypad. Tweets, however, was largely spared from this phenomenon and usually have correct spelling. Among the three media analyzed in this study, Twitter is the only considered medium that is constrained. Conversations in books and movies are supposedly oral conversations that were written down

in the form of a book or a subtitle so their written form should have no effect on them.

We now argue that if conversations are independent of medium, then no significant difference should be observed among conversations in Twitter, books and movies. On the other hand, if differences in a medium is due to an explicit quirk in the medium e.g., an utterance length limit, then conversations in Twitter must be significantly different from conversations in books and movies, but the latter two should not be significantly different from each other. Finally, if conversations are indeed dependent on medium, then conversations in Twitter, movies, and books must be significantly different from each other.

2 Orthographic sentence length and the Brown corpus

The study of sentence lengths in text dates back to the 1939 paper of Udny Yule [16] where it was used to establish authorship. More recently, sentence length has been used to classify text genre by itself [17] or in combination with other text properties [18]. Yule's 1939 paper did not provide the sentence length distribution but several decades after its publication, the distribution was described as log-normal [19, 20, 21] which was later shown by Sichel [22, 23] to be flawed. More recently, Sigurd et al. [24] showed that sentence length distributions may be approximated by a gamma distribution.

In this work, we used the non-standard unit of number of characters (orthographic length), instead of the usual sentence length units of clauses or words, in measuring utterance lengths for ease of comparison with Twitter which has a maximum utterance length in terms of characters. Although the distribution of sentence lengths in terms of words and word lengths in terms of letters can be described by a gamma distribution [24], there is no mathematical guarantee that the distribution of sentence lengths in terms of letters would also follow the same distribution in the general case of different shape and scale parameters of the sentence length (in words) and word length (in letters) distributions. If it can be shown that the sentence length (in letters) distribution can be approximated by a member of the same distribution family as the sentence length (in words), then the use of sentence length comparison using orthographic length is a valid approach.

The Brown corpus [25] consists of about one million words of edited English prose printed during 1961 in the United States [26]. To verify if measuring sentence lengths in terms of characters may be approximated by a gamma distribution, the sentence length (in letters) distribution was simulated, as follows. The word length (in letters) and sentence length (in words) distributions of the tagged Brown corpus was first constructed using the natural language toolkit [27] Python module. In constructing the distributions, only words that contain at least one letter were considered. A gamma distribution given by,

$$\Pr(x) = \frac{x^{\alpha-1} e^{-x/s}}{s^\alpha \Gamma(\alpha)}, \quad (1)$$

where α and s are fitting parameters that describe the shape and ordinate scaling factor, respectively, were then fitted on each distribution using the maximum likelihood estimation [28] feature of the Scipy python module [29]. For each trial, 100,000 sentences were generated following the fitted sentence length (in words) and word length (in letters) distributions. This process was repeated for a total of 100 trials resulting to 100 sentence length in letters histograms. The histograms were converted to a single probability distribution by using the median frequency for each sentence length.

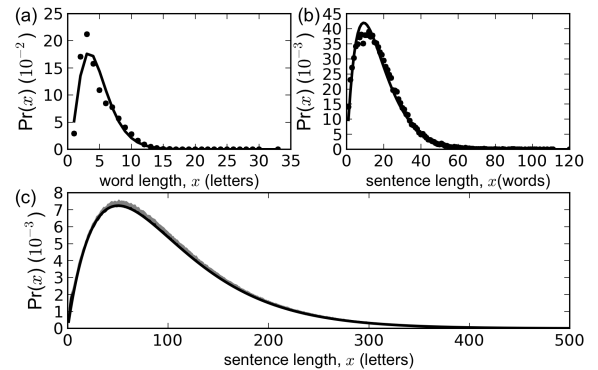


Fig. 1. (a) Word length in letters and (b) sentence length in words distributions of the Brown corpus superimposed with the maximum likelihood estimate of Eq. (1) (solid line). (c) Simulated sentence length (solid dots) in letters distribution using the fitted word length (in letters) and sentence length (in words) distributions of the Brown corpus superimposed with the least-squares fit (solid line) and values within one standard deviation (shaded)

Both the word length (WL, in letters) [Fig. 1(a)] and sentence length (SL, in words) [Fig. 1(c)] distributions of the Brown corpus follow a gamma distribution (WL: $\alpha = 3.43$, $s = 1.39$, $r^2 = 0.948$; SL: $\alpha = 2.09$, $s = 8.44$, $r^2 = 0.989$). The simulated sentence length in letters distribution [Fig. 1(b)] also follows a gamma distribution ($\alpha = 1.98$, $s = 51.2$) but has a much larger s than the sentence length in words which is expected since letters is a smaller syntactic unit than words.

The sentence length distribution in letters thus belongs to the same family of distributions as when measured in words. Since utterance lengths are being compared empirically, the use of orthographic length as a unit of utterance length is therefore valid despite known idiosyncrasies [30] of the English language. Interestingly, the orthographic length was also used by Piantadosi et al. [31] when they showed that word lengths are optimized for efficient communication because it is easier to measure while still being highly correlated with word length in terms of syllables [32].

3 Datasets

Four datasets were used for our analysis: utterances in fictional works in Project Gutenberg (PG)(gutenberg.org), utterances in PG split into sentences (PGS), tweets from Twitter (TWITTER), and utterances in movie subtitles (SUBS) from opensubtitles.org.

PG was generated by extracting utterances—defined as text enclosed in double quotes—from the available works in Project Gutenberg of 50 authors whose selection was roughly based on availability (see Ref. [33] for list of titles, and Ref. [34] for author selection and text parsing details). The resulting dataset consists of about 2.3 million utterances, with zero-length utterances (0.01% of original dataset) removed. The author with the most number of utterances (George Manville Fenn) has 238,640 utterances while the author with the least number of utterances (David Herbert Lawrence) has 1,170 utterances. The median number of utterances is 36,955 utterances per author. When split into sentences, PG is converted to PGS which has about 4.2 million utterances with a median number of utterances equal to 69,311 utterances per author.

Conversations in TWITTER were identified by looking for *replies*, which are Twitter messages (or *tweets*) directed to specific users. We used the convention that *replies* begin with the `@username` of the receiver, e.g., *@bob Hello! How are you?* to filter the tweets for our dataset.¹ Though not in the original design, the use of *replies* emerged as the leading method of addressing a particular person in Twitter [35]. The presence of an `@username` anywhere in the tweet makes that tweet a *mention* [36]. Unlike *mentions*, which appear in the timeline of a user following the sender, a *reply* appears in said user’s timeline only if he follows both sender and receiver of the *reply* message. Thus, conversations are most likely restricted to *replies* to avoid flooding the timeline of people not involved in the discussion. Though *mentions* may carry conversations, we still excluded them from the dataset, as they are more likely non-conversational tweets.

It is possible that a *reply* is not reciprocated, e.g., if it was meant to bring an item, such as a URL, to the attention of another user. This is still considered a conversation because it conveys a short bit of information directly targeted to a certain user. This is similar to someone telling another to “watch out!” or “be careful”: a reply by the other person is not required.

Using the Twitter Streaming application programming interface (API) [37], five one-week sampled public tweets from September 2009 to July 2010 were selected. From the one-week samples composed of around 16.2 million to 57.6 million tweets representing about 15% of public tweets [37], nonzero-length messages were extracted which yielded about 52 million messages or utterances (see Ref. [34] for datasets and parsing details). For better comparison with PG and PGS that have 50 subsets (authors)

each, the weekly datasets were subdivided into ten groups of shuffled hourly data.

SUBS consists of about 14.7 million utterances from 15,809 movies provided by opensubtitles.org. The movie release years span from 1896 to 2010. See Ref. [34] for parsing details and Ref. [38] for the complete list of movies.

4 Utterance length distributions of datasets

Twitter conversations [Fig. 2(a)] have an asymmetric and bimodal utterance length distribution. The left peak (mode) is at 16 characters which we take to be the natural distribution of message lengths i.e., it is the distribution of an unrestricted conversation. Similar to the argument used by Sigurd et al. [24] in their study of word and sentence length distributions of English, Swedish and German texts, and by Cancho and Solé [39] in their work on the origin of Zipf’s law, we posit that the length of an utterance in a conversation is also governed by a trade-off between packing as much information as possible in an utterance and expressing the utterance as quickly as possible: the first objective is biased towards increasing length ($\sim x^{\alpha-1}$) while the other is biased towards decreasing it ($\sim e^{-x}$). Combining the two objectives, the following distribution is obtained: $\sim x^{\alpha-1}e^{-x}$.

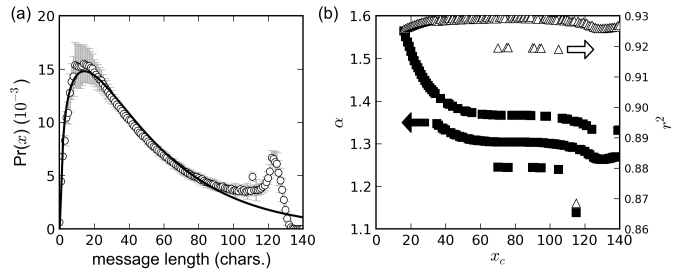


Fig. 2. (a) Message length distribution of sampled tweets with the curve fit having the highest r^2 value ($\alpha = 1.37$, solid line). Error bars are standard deviations from five one-week samples. (b) The α values (filled squares) of the fit from $x = 0$ to x_c using Eq. (2) and its corresponding r^2 (unfilled triangles).

To account for a strict length limit for Twitter messages, the natural utterance length distribution was estimated by fitting a more general equation using a modified Levenberg-Marquardt least squares algorithm [29] to the utterance length distribution from $x = 0$ to a cut-off length $x_c \in [16, 140]$ [Fig. 2(b)],

$$\Pr(x) = \frac{\tilde{x}^{\alpha-1} e^{-\tilde{x}}}{\Gamma(\alpha)}, \quad (2)$$

where $\tilde{x} = (x - x_0)/s$ is the scaled utterance length x , while α , x_0 and s are fitting parameters that describe the shape, translation and ordinate scaling factor, respectively. This method of estimation assumes that the mixing parameter of the bimodal distribution is almost one in favor of the natural utterance length distribution. A bimodal distribution fitted using expectation maximization was not

¹ The current Twitter API supports a method for explicitly classifying a tweet as a *reply* but this was not yet widely available and followed when our data were gathered.

utilized because of a lack of an explicit model of the truncation distribution. Our goal is to estimate the median of the natural utterance length distribution so a resulting non-normalized unimodal distribution is acceptable.

When α approaches one, Eq. (2) approaches an exponential distribution. The range of acceptable values of $\alpha \in [1.1, 1.6]$, [$r^2 \in (0.86, 0.93)$] for the Twitter dataset corresponds to a 57-order-of-magnitude increase in likelihood of finding an utterance length of $x = x_c = 140$ chars. compared to an exponentially decaying curve in the absence of a Twitter-imposed limit (see Ref. [27] for the fitting parameters distributions). However, another peak was found at 124 characters due to the 140-character limit, a limit that is absent in the other datasets, and is attributed to various tweet-shortening schemes. The absence of a length limit results to unimodal utterance length distributions for PG, PGS and SUBS [Fig. 3].

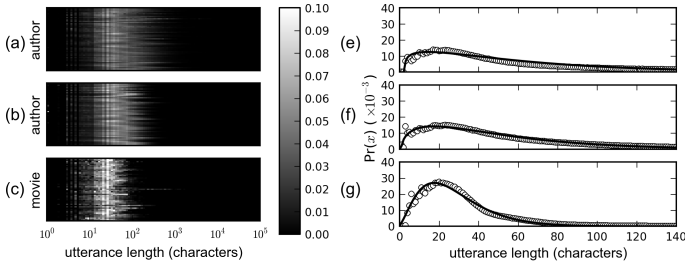


Fig. 3. Utterance length distributions of (a) different authors in PG (b) different authors in PGS and (c) 50 randomly selected movies in SUBS. Distribution of utterance lengths over the entire (d) PG, (e) PGS and (f) SUBS datasets fitted with Eq. (2).

Conversations in movies (interquartile range IQR = difference between the 3rd and 1st quartiles = 21 chars.) are of more uniform length than those in books (PG IQR median = 88 chars, PGS IQR median = 50 chars.). The much smaller SUBS IQR median compared to that of TWITTER (IQR median = 46 chars.) or that of its best fit of Eq. (2) (IQR median = 50 chars.) suggests that conversations in movies are less dependent on author style while the much larger IQR medians of PG and PGS point to a stronger dependence of these media on author style.

To minimize the effect of unequal author or movie utterances, and of noise due to differences in spelling and punctuation, Eq. (2) was fitted to PG, PGS and SUBS by computing for the normalized histogram of each author or movie then using the average probability for each utterance length as the probability density function to be fitted using least squares. Based on the fit of Eq. (2) ($\alpha = 1.48$, $x_0 = 0.862$, $s = 34.4$, $r = 0.984$), the PGS utterance length distribution [Fig. 3(e)] seems to be a horizontally compressed TWITTER best fit curve ($\alpha = 1.37$, $x_0 = 0.86$, $s = 36.4$) because of a smaller s value. The PG utterance length distribution has a fatter tail [$1 - F(140) = 0.0896$; Fig. 3(d)] than that of the PGS utterance length distribution [$1 - F(140) = 0.0427$], and only its tail fits Eq. (2) quite well ($\alpha = 1.24$, $x_0 = 2.63$, $s = 48.6$, $r^2 = 0.970$). In contrast, the entire SUBS median distribution [$\alpha = 2.71$,

$x_0 = 0.87$, $s = 10.7$, $r^2 = 0.988$; Fig. 2(f)] fits Eq. (2) and has almost no tail ($1 - F(140) = 1.19 \times 10^{-4}$). Thus, all datasets share the same distribution family as the Brown sentence length in words distribution further giving credence to the validity of the use of characters as a unit of utterance length.

The mean length of utterance (MLU) is used to evaluate the level of language development of a child [40,41]. However, the use of the mean as a measure of central tendency is invalid because the utterance length distribution is very skewed to the right. The mode of a gamma distribution [Eq. (2)] is given by $(\alpha - 1)s + x_0$ but it does not appear to be correlated with s [Fig. 4(a)]. In contrast, the median, though not having a closed form equation for a gamma distribution, appears to be more correlated with s [Fig. 4(b)]: a larger median roughly implies a larger spread. The median, therefore, allows us to simultaneously describe both the location and scale of the utterance length distribution.

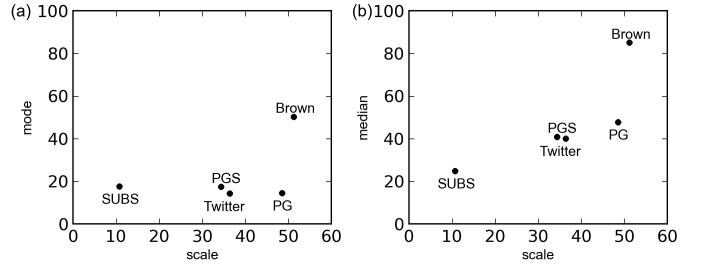


Fig. 4. Mode and median of the distribution fits. (a) Mode and (b) median of the fit of each distribution plotted against s .

For the rest of this paper, the median utterance length and its median were used to describe each utterance length distribution. These measures are suitable for comparison between datasets because both are insensitive to outliers (robust) and do not assume a distribution (nonparametric). Any author dependence or deviation from a gamma distribution of the data would therefore not affect the results [34]. Tests for significant differences were performed using the Mann-Whitney U test [42] with continuity correction because the distributions being compared are discrete and skewed.

5 Utterance length and sample size

TWITTER, PGS and SUBS were subsampled (with replacement) such that the sample size would be the same for each author's sample size in PG. By taking the distribution of subsample medians (Fig. 5) which is analogous to taking the distribution of sample means from normally-distributed data, we found that the median median utterance length (analogous to mean of sample means) of SUBS (25 chars.) is very different from that of TWITTER (38 chars.), PG (48 chars.) and PGS (41 chars.).

Notably, the median median utterance length value of SUBS of 25 chars., which is not related to the existing max-

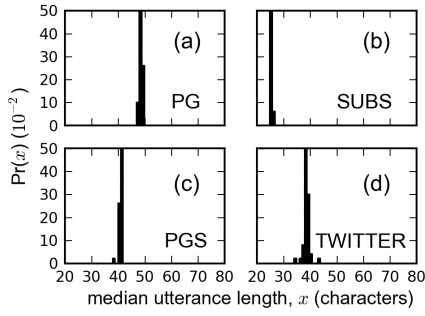


Fig. 5. Distribution of median utterance lengths (median median utterance length: dashed lines) for (a) PG, (b) SUBS, (c) PGS and (d) TWITTER. The median utterance length data in (d) was estimated from the natural utterance length distribution of each TWITTER subset.

imum subtitle line length of 32-34 characters (Ofcom regulation [43]), points to a fundamental difference in how the verbal medium is used in movies.

The median utterance length distribution of all datasets are significantly different from each other (see Ref. [27] for complete test results between each pair of dataset). Since the PGS median distribution is significantly different from the SUBS median distribution, conversational sentences in books are not the same as conversational sentences in movies though we posit that conversations in movies are closer to that of actual transcribed speech. TWITTER utterance lengths are stochastically smaller than PG and PGS but differ significantly from SUBS suggesting that Twitter is a less formal medium. We surmise that the smaller length is due to the more spontaneous and less formal tone of Twitter conversations than those in books.

To investigate the effect of sample size N on the median utterance length, each dataset was sampled (with replacement) into 50 groups each having N utterances. Similar to word frequency distributions that are dependent on N [44], the spread in, but not the location of, the medians distribution decreases as N increases (Fig. 6) for all datasets. At $N = 10^5$ utterances, the median value of SUBS collapsed to a single value of 25 characters. At $N = 10^6$ utterances, PG and PGS collapsed to different single median utterance length values of 48 and 41 characters, respectively, while TWITTER falls into two unique values of 38 and 39 characters.

The median utterance length distribution of SUBS is very different from the median utterance length distribution of the other datasets—it can be clearly distinguished from them even if the sample size is only $N = 100$ utterances (Fig. 6). PG and PGS median utterance length distributions are already distinguishable from each other but both overlap with TWITTER at $N = 100$ utterances. The PG, PGS and TWITTER median utterance length distributions do not overlap only at $N = 10^4$ utterances, thus giving us the required minimum sample size for meaningful comparison across communication media as a function of time (see Ref. [34] for complete test results).

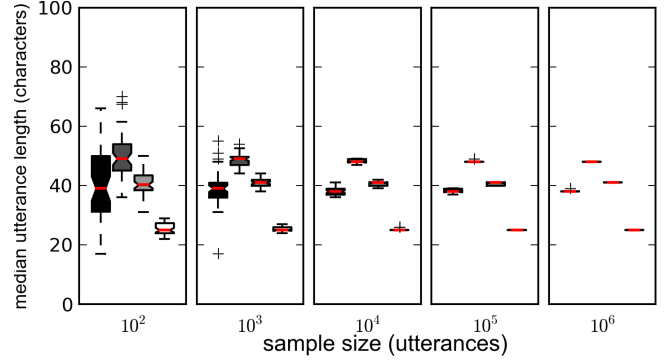


Fig. 6. Distribution of median utterance length in subsampled TWITTER (black), PG (dark gray), PGS (light gray) and SUBS (unfilled).

6 Utterance length through time

The median median utterance length in both PG [Fig. 7(a)] (slope = -0.266 chars./yr, $r^2 = 0.903$, $p < 10^{-3}$ two-sided) and PGS [Fig. 7(b)] (slope = -0.189 chars./yr, $r^2 = 0.814$, $p < 10^{-3}$ two-sided) decreases with time but is not correlated with size (PG Spearman $\rho^2 < 10^{-3}$; PGS Spearman $\rho^2 = 0.00524$).

On the other hand, the median utterance length of SUBS [Fig. 7(c)] remains almost constant (~ 27 chars.) in time (slope = -1.897×10^{-3} chars./yr, $r^2 = 0.121$, $p < 10^{-3}$ two-sided) except for a conspicuous rise and increased spread in the median utterance length at around 1920 that does not flatten out even if the window size is increased from 1 year to 5 years [Fig. 7(d)]. The bump is likely due to the availability of “talking pictures” and commercial television starting in the late 1920s. The silent movies prior to their release have a different “conversation signature” from those of “talkies”.

The temporal behavior of TWITTER was not studied because TWITTER spans only a few weeks.

7 Conclusion

Though we do not usually notice the medium-dependence of conversations, we showed that conversations, as measured by orthographic utterance length, are slowly shortening in time within media but are drastically different across different media. These are fundamental differences that are effects not just of the milieu, but of the medium itself. Evolving technologies that lead to changes in communication media seemingly lead us to adapt our conversations, rather than such a technology suffering an early demise because it cannot adapt to our natural use of language. An extreme case in point is the short message service (SMS) or “texting.” Originally designed with a character limit of 160 such that most sentences would fit in a single text message [45], but with an “access a letter via numerical keypad” constraint—it became a popular form of communication [46] with its own lingo [47]. Clearly, adaptation occurs with changing medium and sometimes with unexpected side-effects.

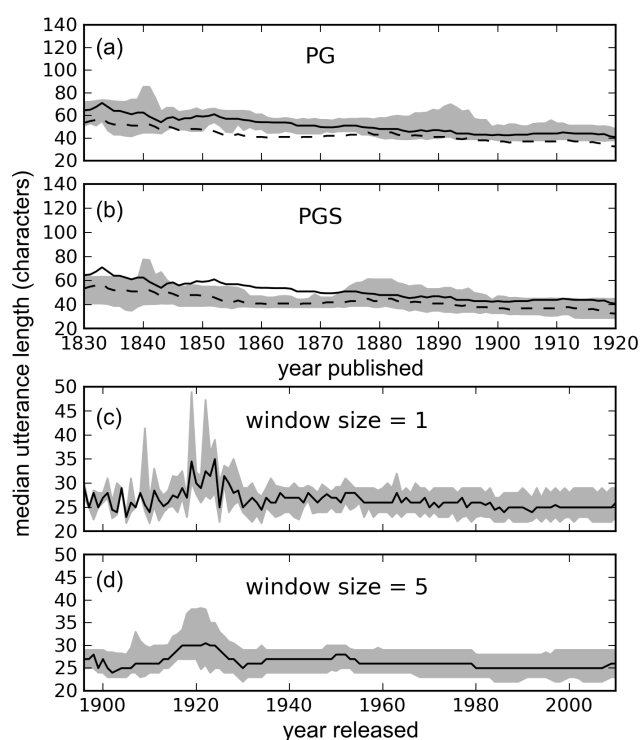


Fig. 7. Median utterance length distribution of (a) PG and (b) PGS with window size of 10 years, and SUBS with window size of (c) 1 year and (d) 5 years. Only books with at least 1,000 utterances were considered. Publication years were retrieved from the US Library of Congress. The window sizes were selected so that the plots do not change appreciably when the window size is varied slightly. First to third quartiles (shaded), PG median utterance length (a-b, solid line), PGS median utterance length (a-b, dashed line).

We thank the administrator of opensubtitles.org for providing us the text version of their English-languages movies subtitles. This work is supported by a grant from the UP Diliman Office of the Vice Chancellor for Research and Development and by an Amazon AWS Education grant.

References

1. R. Goulden, P. Nation, J. Read, *Appl. Linguistics*, **11**, 341 (1990).
2. P. Nation, R. Waring, in *Vocabulary: Description, Acquisition and Pedagogy*, edited by N. Schmitt, M. McCarthy (Cambridge University Press, Cambridge, 1997).
3. C. Browne, G. Cih, B. Culligan, "Measuring vocabulary size via online technology" (2007), <http://www.lexxica.com> [Retrieved 08-12-2012]
4. D. P. Hayes, M. G. Ahrens, *J. of Child Lang.*, **15**, 395 (1988).
5. S. Hill, N. Launder, *Australian J. of Lang. and Lit.*, **33**, 240 (2010).
6. W. Chafe, D. Tannen, *Annual Rev. of Anthropology*, **16**, 383 (1987).
7. R. Wooffitt, *Conversation analysis and discourse analysis: A comparative and critical introduction* (Sage Publications Ltd, 2005).
8. W. Sack, *J. Manage. Inf. Syst.*, **17**, 73 (2000).
9. A. Warner, *Complementation in Middle English and the Methodology of Historical Syntax: A Study of the Wyclifite Sermons* (Taylor & Francis, 1982).
10. A. Warner, *Language Variation and Change*, **17**, 257 (2005).
11. D. Lorenz, in *ICAME 33: Corpora at the centre and crossroads of English linguistics, Leuven, 2012* (University of Leuven, 2012), p. 185.
12. M. Kytö and T. Walker, *Guide to A Corpus of English Dialogues, 1560-1760*, (Uppsala Universitet, 2006).
13. M. Davies, *The Corpus of Historical American English: 400 million words, 1810-2009*, (2010), <http://corpus.byu.edu/coha/> [Retrieved 23-07-2012]
14. J. A. Smith, C. Kelly, *Computers and the Humanities*, **36**, 411 (2002).
15. D. Crystal, *Language and the Internet*, 2nd edn. (Cambridge Univ Press, 2006).
16. G. U. Yule, *Biometrika*, **30**, 363 (1939).
17. E. Kelih, P. Grzybek, G. Antic, E. Stadlobör, in *From Data and Information Analysis to Knowledge Engineering*, edited by M. Spiliopoulou, R. Kruse, C. Borgelt, A. Nurnberger, W. Gaul (Springer Berlin Heidelberg, 2006).
18. T. Copeck, K. Barker, S. Delisle, S. Szpakowicz, in *TALN-2000: Actes de la 7e Conference Annuelle sur le Traitement Automatique des Langues Naturelle, Laussane, 2000*.
19. C. B. Williams, *Biometrika*, **31**, 356 (1940).
20. C. B. Williams, *Style and vocabulary: numerical studies* (Griffin, 1970).
21. W. C. Wake, *J. Royal Statistical Soc. A*, **120**, 331 (1957).
22. H. S. Sichel, *J. Royal Statistical Soc. A*, **137**, 25 (1974).
23. P. Grzybek, in *Contributions to the Science of Text and Language*, edited by P. Grzybek (Springer Netherlands, Dordrecht, 2006).
24. B. Sigurd, M. Eeg-Olofsson, J. van de Weijer, *Studia Linguistica*, **58**, 37 (2004).
25. H. Kučera and W. Francis, *Computational analysis of present-day American English* (Dartmouth Publishing Group, 1967).
26. W. Francis, H. Kucera, *Brown Corpus Manual of Information* (Brown University, 1979), <http://icame.uib.no/brown/bcm.html> [Retrieved 19-01-2012]
27. S. Bird, E. Loper, E. Klein, *Natural Language Processing with Python* (O'Reilly Media Inc., 2009).
28. L. Wasserman, *All of Statistics: A Concise Course in Statistical Inference* (Springer, 2004).
29. E. Jones, T. Oliphant, P. Peterson, *et al.*, "SciPy: Open Source Scientific Tools for Python," 2001, <http://www.scipy.org/> [Retrieved 19-04-2011]
30. "English spelling: You write potato, i write ghoughpteighbteau," *The Economist* (2008).
31. S. T. Piantadosi, H. Tily, E. Gibson, *Proc. Natl. Acad. Sci. U. S. A.*, **108**, 3526 (2011).
32. U. Strauss, P. Grzybek, G. Altmann, in *Contributions to the Science of Text and Language*, edited by P. Grzybek, (Springer-Verlag, Berlin/Heidelberg, 2006).
33. C. M. Alis, M. T. Lim, "Supplementary material: PG authors list." (2012), <http://www.nip.upd.edu.ph/ipd/data/conversations/pg-authorslist.csv>

34. C. M. Alis, M. T. Lim, "Supplementary material: Adaptation of fictional and online conversations to communication media" (2012), http://www.nip.upd.edu.ph/ipl/data/conversations/epjb_si.pdf [Retrieved 27-09-2012]
35. E. Williams, "How @replies work on twitter (and how they might)" (2008), <http://blog.twitter.com/2008/05/how-replies-work-on-twitter-and-how.html> [Retrieved 25-09-2012]
36. Twitter Help Center, "What are @replies and mentions?" (2012), <https://support.twitter.com/articles/14023-what-are-replies-and-mentions> [Retrieved 25-09-2012]
37. J. Kalucki, "Streaming API documentation," (2010), <http://apiwiki.twitter.com/w/page/22554673/Streaming-API-Documentation?rev=1268351420> [Retrieved 15-04-2012]
38. C. M. Alis, M. T. Lim, "Supplementary material: SUBS movie list" (2012), <http://www.nip.upd.edu.ph/ipl/data/conversations/subs-movielist.csv> [Retrieved 27-09-2012]
39. R. F. i. Cancho, R. V. Solé, *Proc. Natl. Acad. Sci. U. S. A.*, **100** 788 (2003)
40. T. Klee, M. D. Fitzgerald, *J. of Child Lang.*, **12** 251 (1985)
41. C. A. Dollaghan, T. F. Campbell, J. L. Paradise, H. M. Feldman, J. E. Janosky, D. N. Pitcairn, M. Kurs-Lasky, *J. Speech Lang. Hear. Res.*, **42** 1432 (1999)
42. H. B. Mann, D. R. Whitney, *Ann. Math. Stat.*, **18**, 50 (1947)
43. Independent Television Commission, "ITC Guidance on standards for subtitling," (1999), http://www.ofcom.org.uk/static/archive/itc/itc_publications/codes_guidance/index.asp.html [Retrieved 15-04-2012]
44. S. Bernhardsson, L. E. C. da Rocha, P. Minnhagen, *New J. of Phys.*, **11**, 123015 (2009)
45. M. Milian, "Why text messages are limited to 160 characters," *Los Angeles Times*, May 2009.
46. "ictDATA.org: top SMS 2009." (2010), <http://www.ictdata.org/2010/10/top-sms-2009.html> [Retrieved 21-02-2012]
47. C. Thurlow, *Discourse Analysis Online*, **1**(1) (2003), <http://www.shu.ac.uk/daol/articles/v1/n1/a3/thurLOW2002003-paper.html> [Retrieved 21-02-2012]